



TITLE:

Asymptotic properties of support vector machines in HDLSS settings (Bayes Inference and Its Related Topics)

AUTHOR(S):

Nakayama, Yugo; Yata, Kazuyoshi; Aoshima, Makoto

CITATION:

Nakayama, Yugo ...[et al]. Asymptotic properties of support vector machines in HDLSS settings (Bayes Inference and Its Related Topics). 数理解析研究所講究録 2017, 2047: 10-18

ISSUE DATE:

2017-10

URL:

<http://hdl.handle.net/2433/237025>

RIGHT:

Asymptotic properties of support vector machines in HDLSS settings

Yugo Nakayama

Graduate School of Pure and Applied Sciences
University of Tsukuba

Kazuyoshi Yata

Institute of Mathematics
University of Tsukuba

Makoto Aoshima

Institute of Mathematics
University of Tsukuba

Abstract

In this paper, we consider asymptotic properties of the support vector machine (SVM) in high-dimension, low-sample-size (HDLSS) settings. We first show that the linear SVM holds a consistency property in which misclassification rates tend to zero as the dimension goes to infinity under certain severe conditions. Next, we consider a non-linear SVM based on the Gaussian kernel in HDLSS settings. We also show that the non-linear SVM holds the consistency property under mild conditions. Finally, we check the performance of the SVMs by numerical simulations.

Keywords and phrases: Hard-margin SVM; Large p small n ; Radial basis function kernel

1 Introduction

High-dimension, low-sample-size (HDLSS) data situations occur in many areas of modern science such as genetic microarrays, medical imaging, text recognition, finance, chemometrics, and so on. Suppose we have independent and d -variate two populations, π_i , $i = 1, 2$, having an unknown mean vector μ_i and unknown covariance matrix Σ_i ($\geq O$). We assume that $\text{tr}(\Sigma_i)/d \in (0, \infty)$ as $d \rightarrow \infty$ for $i = 1, 2$. Here, for a function, $f(\cdot)$, “ $f(d) \in (0, \infty)$ as $d \rightarrow \infty$ ” implies $\liminf_{d \rightarrow \infty} f(d) > 0$ and $\limsup_{d \rightarrow \infty} f(d) < \infty$. Let $\Delta = \|\mu_1 - \mu_2\|^2$, where $\|\cdot\|$ denotes the Euclidean norm. We assume that $\limsup_{d \rightarrow \infty} \Delta/d < \infty$. We have independent and identically distributed (i.i.d.) observations, $\mathbf{x}_{i1}, \dots, \mathbf{x}_{in_i}$, from each π_i . We

assume $n_i \geq 2$, $i = 1, 2$. Let \mathbf{x}_0 be an observation vector of an individual belonging to one of the two populations. We assume \mathbf{x}_0 and \mathbf{x}_{ij} s are independent. Let $N = n_1 + n_2$.

In the HDLSS context, Hall et al. [6] and Marron et al. [7] considered distance weighted classifiers. Aoshima and Yata [2] and Chan and Hall [5] considered distance-based classifiers. In particular, Aoshima and Yata [2] gave the misclassification rate adjusted classifier for multiclass, high-dimensional data in which misclassification rates are no more than specified thresholds. On the other hand, Aoshima and Yata [1, 3] considered geometric classifiers based on a geometric representation of HDLSS data. Aoshima and Yata [4] considered quadratic classifiers in general and discussed asymptotic properties and optimality of the classifiers under high-dimension, non-sparse settings. They showed that the misclassification rates tend to 0 as d increases, i.e.,

$$e(i) \rightarrow 0 \text{ as } d \rightarrow \infty \text{ for } i = 1, 2 \quad (1)$$

under the non-sparsity such as $\Delta \rightarrow \infty$ as $d \rightarrow \infty$, where $e(i)$ denotes the error rate of misclassifying an individual from π_i into the other class. We call (1) “the consistency property”.

In the field of machine learning, there are many studies about the classification in the context of supervised learning. A typical method is the support vector machine (SVM). The SVM has versatility and effectiveness both for low-dimensional and high-dimensional data. See Schölkopf and Smola [9] and Vapnik [10] for the details. Even though the SVM is quite popular, its asymptotic properties seem to have not been studied sufficiently. Recently, Nakayama et al. [8] investigated asymptotic properties of a linear SVM for HDLSS data.

In this paper, we investigate linear and non-linear SVMs in the HDLSS context where $d \rightarrow \infty$ while N is fixed. In Section 2, we show that the linear SVM holds (1) under certain severe conditions. In Section 3, we consider a non-linear SVM based on the Gaussian kernel in HDLSS settings. We also show that the non-linear SVM holds (1) under mild conditions. Finally, we check the performance of the SVMs by numerical simulations.

2 Linear SVM in HDLSS settings

In this section, we give asymptotic properties of the linear SVM in HDLSS settings. Since HDLSS data are linearly separable by a hyperplane, we consider the hard-margin linear SVM.

2.1 Hard-margin linear SVM

We consider the following linear classifier:

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b, \quad (2)$$

where \mathbf{w} is a weight vector and b is an intercept term. Let us write that $(\mathbf{x}_1, \dots, \mathbf{x}_N) = (\mathbf{x}_{11}, \dots, \mathbf{x}_{1n_1}, \mathbf{x}_{21}, \dots, \mathbf{x}_{2n_2})$. Let $t_j = -1$ for $j = 1, \dots, n_1$ and $t_j = 1$ for $j = n_1 + 1, \dots, N$. The hard-margin SVM is defined by maximizing the smallest distance of all observations to the

separating hyperplane. The optimization problem of the SVM can be written as follows:

$$\underset{\mathbf{w}, b}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{subject to} \quad t_j(\mathbf{w}^T \mathbf{x}_j + b) \geq 1, j = 1, \dots, N.$$

A Lagrangian formulation is given by

$$L(\mathbf{w}, b; \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{j=1}^N \alpha_j \{t_j(\mathbf{w}^T \mathbf{x}_j + b) - 1\},$$

where $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_N)^T$ and α_j s are Lagrange multipliers. By differentiating the Lagrangian formulation with respect to \mathbf{w} and b , we obtain the following conditions:

$$\mathbf{w} = \sum_{j=1}^N \alpha_j t_j \mathbf{x}_j \quad \text{and} \quad \sum_{j=1}^N \alpha_j t_j = 0.$$

After substituting them into $L(\mathbf{w}, b; \boldsymbol{\alpha})$, we obtain the dual form:

$$L(\boldsymbol{\alpha}) = \sum_{j=1}^N \alpha_j - \frac{1}{2} \sum_{j=1}^N \sum_{k=1}^N \alpha_j \alpha_k t_j t_k \mathbf{x}_j^T \mathbf{x}_k.$$

The optimization problem can be transformed into the following:

$$\underset{\boldsymbol{\alpha}}{\operatorname{argmax}} L(\boldsymbol{\alpha})$$

subject to

$$\alpha_j \geq 0, j = 1, \dots, N, \quad \text{and} \quad \sum_{j=1}^N \alpha_j t_j = 0. \quad (3)$$

Let us write that

$$\hat{\boldsymbol{\alpha}} = (\hat{\alpha}_1, \dots, \hat{\alpha}_N)^T = \underset{\boldsymbol{\alpha}}{\operatorname{argmax}} L(\boldsymbol{\alpha}) \quad \text{subject to (3)}.$$

There exist some \mathbf{x}_j s satisfying that $t_j y(\mathbf{x}_j) = 1$ (i.e., $\hat{\alpha}_j \neq 0$). Such \mathbf{x}_j s are called the support vector. Let $\hat{S} = \{j | \hat{\alpha}_j \neq 0, j = 1, \dots, N\}$ and $N_{\hat{S}} = \#\hat{S}$, where $\#A$ denotes the number of elements in a set A . The intercept term is given by

$$\hat{b} = \frac{1}{N_{\hat{S}}} \sum_{j \in \hat{S}} \left(t_j - \sum_{k \in \hat{S}} \hat{\alpha}_k t_k \mathbf{x}_j^T \mathbf{x}_k \right).$$

Then, the linear classifier in (2) is defined by

$$\hat{y}(\mathbf{x}) = \sum_{k \in \hat{S}} \hat{\alpha}_k t_k \mathbf{x}_k^T \mathbf{x} + \hat{b}. \quad (4)$$

Finally, in the SVM, one classifies \mathbf{x}_0 into π_1 if $\hat{y}(\mathbf{x}_0) < 0$ and into π_2 otherwise. See Vapnik [10] for the details.

2.2 Asymptotic properties of the linear SVM in the HDLSS context

We assume the following assumptions:

$$(A-i) \quad \frac{\text{Var}(\|\mathbf{x}_{ik} - \boldsymbol{\mu}_i\|^2)}{\Delta^2} \rightarrow 0 \text{ as } d \rightarrow \infty \text{ for } i = 1, 2;$$

$$(A-ii) \quad \frac{\text{tr}(\boldsymbol{\Sigma}_i^2)}{\Delta^2} \rightarrow 0 \text{ as } d \rightarrow \infty \text{ for } i = 1, 2.$$

Note that $\text{Var}(\|\mathbf{x}_{ik} - \boldsymbol{\mu}_i\|^2) = 2\text{tr}(\boldsymbol{\Sigma}_i^2)$ when π_i is Gaussian, so that (A-i) and (A-ii) are equivalent when π_i s are Gaussian. Let

$$\delta = \Delta + \frac{\text{tr}(\boldsymbol{\Sigma}_1)}{n_1} + \frac{\text{tr}(\boldsymbol{\Sigma}_2)}{n_2} \quad \text{and} \quad \kappa = \frac{\text{tr}(\boldsymbol{\Sigma}_1)}{n_1} - \frac{\text{tr}(\boldsymbol{\Sigma}_2)}{n_2}.$$

Then, Nakayama et al. [8] gave the following results.

Lemma 2.1 ([8]). *Under (A-i) and (A-ii), it holds that as $d \rightarrow \infty$*

$$\begin{aligned} \hat{\alpha}_j &= \frac{2}{\delta n_1} \{1 + o_p(1)\} \quad \text{for } j = 1, \dots, n_1; \quad \text{and} \\ \hat{\alpha}_j &= \frac{2}{\delta n_2} \{1 + o_p(1)\} \quad \text{for } j = n_1 + 1, \dots, N. \end{aligned}$$

Furthermore, it holds that as $d \rightarrow \infty$

$$\hat{y}(\mathbf{x}_0) = \frac{(-1)^i \Delta}{\delta} + \frac{\kappa}{\delta} + o_p\left(\frac{\Delta}{\delta}\right) \quad \text{when } \mathbf{x}_0 \in \pi_i, i = 1, 2.$$

From Lemma 2.1, it holds that as $d \rightarrow \infty$

$$\frac{\delta}{\Delta} \hat{y}(\mathbf{x}_0) = (-1)^i + \frac{\kappa}{\Delta} + o_p(1) \tag{5}$$

when $\mathbf{x}_0 \in \pi_i, i = 1, 2$. Hence, “ κ/Δ ” is the bias term of the (normalized) SVM. We consider the following assumption:

$$(A-iii) \quad \limsup_{d \rightarrow \infty} \frac{|\kappa|}{\Delta} < 1.$$

Then, Nakayama et al. [8] gave the following results.

Theorem 2.1 ([8]). *Under (A-i) to (A-iii), the linear SVM holds (1).*

Corollary 2.1 ([8]). *Under (A-i) and (A-ii), the linear SVM holds the following properties:*

$$\begin{aligned} e(1) &\rightarrow 1 \quad \text{and} \quad e(2) \rightarrow 0 \quad \text{as } d \rightarrow \infty \quad \text{if} \quad \liminf_{d \rightarrow \infty} \frac{\kappa}{\Delta} > 1; \quad \text{and} \\ e(1) &\rightarrow 0 \quad \text{and} \quad e(2) \rightarrow 1 \quad \text{as } d \rightarrow \infty \quad \text{if} \quad \limsup_{d \rightarrow \infty} \frac{\kappa}{\Delta} < -1. \end{aligned}$$

We expect from (5) that, for sufficiently large d , $e(1)$ and $e(2)$ for the SVM become small and $e(1)$ (or $e(2)$) is larger than $e(2)$ (or $e(1)$) if $\kappa/\Delta > 0$ (or $\kappa/\Delta < 0$). In addition, from Corollary 2.1, if $\liminf_{d \rightarrow \infty} |\kappa|/\Delta > 1$, one should not use the SVM. In order to overcome the difficulties, Nakayama et al. [8] proposed a bias-corrected SVM (BC-SVM). They showed that the BC-SVM gives preferable performances even when (A-iii) is not met.

3 Non-linear SVM in HDLSS settings

In this section, we consider a non-linear SVM based on the Gaussian kernel. We give asymptotic properties of the non-linear SVM in HDLSS settings.

The optimization problem of the non-linear SVM can be written as follows: Let

$$L_*(\alpha) = \sum_{j=1}^N \alpha_j - \frac{1}{2} \sum_{j=1}^N \sum_{k=1}^N \alpha_j \alpha_k t_j t_k \exp \left(- \frac{\|\mathbf{x}_j - \mathbf{x}_k\|^2}{\gamma} \right),$$

where $\gamma > 0$ is a tuning parameter. The optimization problem can be transformed into the following:

$$\operatorname{argmax}_{\alpha} L_*(\alpha)$$

subject to (3). Let us write that

$$\tilde{\alpha} = (\tilde{\alpha}_1, \dots, \tilde{\alpha}_N)^T = \operatorname{argmax}_{\alpha} L_*(\alpha) \text{ subject to (3).}$$

Let $\tilde{S} = \{j | \tilde{\alpha}_j \neq 0, j = 1, \dots, N\}$ and $N_{\tilde{S}} = \#\tilde{S}$. The intercept term is given by

$$\tilde{b} = \frac{1}{N_{\tilde{S}}} \sum_{j \in \tilde{S}} \left(t_j - \sum_{k \in \tilde{S}} \tilde{\alpha}_k t_k \exp \left(- \frac{\|\mathbf{x}_j - \mathbf{x}_k\|^2}{\gamma} \right) \right).$$

Then, the classifier is given by

$$\tilde{y}(\mathbf{x}) = \sum_{k \in \tilde{S}} \tilde{\alpha}_k t_k \exp \left(- \frac{\|\mathbf{x}_k - \mathbf{x}\|^2}{\gamma} \right) + \tilde{b}. \quad (6)$$

Finally, in the non-linear SVM, one classifies \mathbf{x}_0 into π_1 if $\tilde{y}(\mathbf{x}_0) < 0$ and into π_2 otherwise.

We assume the following condition for γ :

(A-iv) $\gamma/d \in (0, \infty)$ as $d \rightarrow \infty$.

Let

$$c_i = \exp \left(- \frac{2\operatorname{tr}(\Sigma_i)}{\gamma} \right), \quad i = 1, 2; \quad \text{and} \quad c_3 = \exp \left(- \frac{\operatorname{tr}(\Sigma_1) + \operatorname{tr}(\Sigma_2) + \Delta}{\gamma} \right).$$

Let $\Delta_* = c_1 + c_2 - 2c_3$, $\delta_* = \Delta_* + \sum_{i=1}^2 (1 - c_i)/n_i$ and $\kappa_* = (1 - c_1)/n_1 - (1 - c_2)/n_2$. Here, we assume the following assumptions:

$$(A-v) \quad \frac{\text{Var}(\|\mathbf{x}_{ij} - \boldsymbol{\mu}_i\|^2)}{d^2 \Delta_*^2} \rightarrow 0 \text{ as } d \rightarrow \infty \text{ for } i = 1, 2;$$

$$(A-vi) \quad \frac{\text{tr}(\boldsymbol{\Sigma}_i^2)}{d^2 \Delta_*^2} \rightarrow 0 \text{ as } d \rightarrow \infty \text{ for } i = 1, 2.$$

We have the following result.

Lemma 3.1. *Assume (A-iv) to (A-vi). It holds that as $d \rightarrow \infty$*

$$\begin{aligned} \tilde{\alpha}_j &= \frac{2}{\delta_* n_1} \{1 + o_p(1)\} \quad \text{for } j = 1, \dots, n_1; \quad \text{and} \\ \tilde{\alpha}_j &= \frac{2}{\delta_* n_2} \{1 + o_p(1)\} \quad \text{for } j = n_1 + 1, \dots, N. \end{aligned}$$

Furthermore, it holds that as $d \rightarrow \infty$

$$\tilde{y}(\mathbf{x}_0) = \frac{(-1)^i \Delta_*}{\delta_*} + \frac{\kappa_*}{\delta_*} + o_p\left(\frac{\Delta_*}{\delta_*}\right) \quad \text{when } \mathbf{x}_0 \in \pi_i \text{ for } i = 1, 2. \quad (7)$$

We consider the following assumption:

$$(A-vii) \quad \limsup_{d \rightarrow \infty} \frac{|\kappa_*|}{\Delta_*} < 1.$$

Then, from Lemma 3.1, we have the following result.

Theorem 3.1. *Under (A-iv) to (A-vii), the non-linear SVM holds (1).*

Now, we consider the following conditions:

$$\text{Var}(\|\mathbf{x}_{ij} - \boldsymbol{\mu}_i\|^2) = O\{\text{tr}(\boldsymbol{\Sigma}_i^2)\} \text{ and } \text{tr}(\boldsymbol{\Sigma}_i^2)/d^2 \rightarrow 0 \text{ as } d \rightarrow \infty \text{ for } i = 1, 2. \quad (8)$$

We note that

$$\Delta_* \geq [\exp\{-\text{tr}(\boldsymbol{\Sigma}_1)/\gamma\} - \exp\{-\text{tr}(\boldsymbol{\Sigma}_2)/\gamma\}]^2.$$

If one can assume that $\liminf_{d \rightarrow \infty} |\text{tr}(\boldsymbol{\Sigma}_1)/\text{tr}(\boldsymbol{\Sigma}_2) - 1| > 0$, it follows $\liminf_{d \rightarrow \infty} \Delta_* > 0$ under (A-iv), so that (A-v) and (A-vi) hold under (8). Thus the non-linear SVM has the consistency even when $\boldsymbol{\mu}_1 = \boldsymbol{\mu}_2$. We emphasize that the non-linear SVM based on the Gaussian kernel draws information about heteroscedasticity via the difference of $\text{tr}(\boldsymbol{\Sigma}_i)$ s.

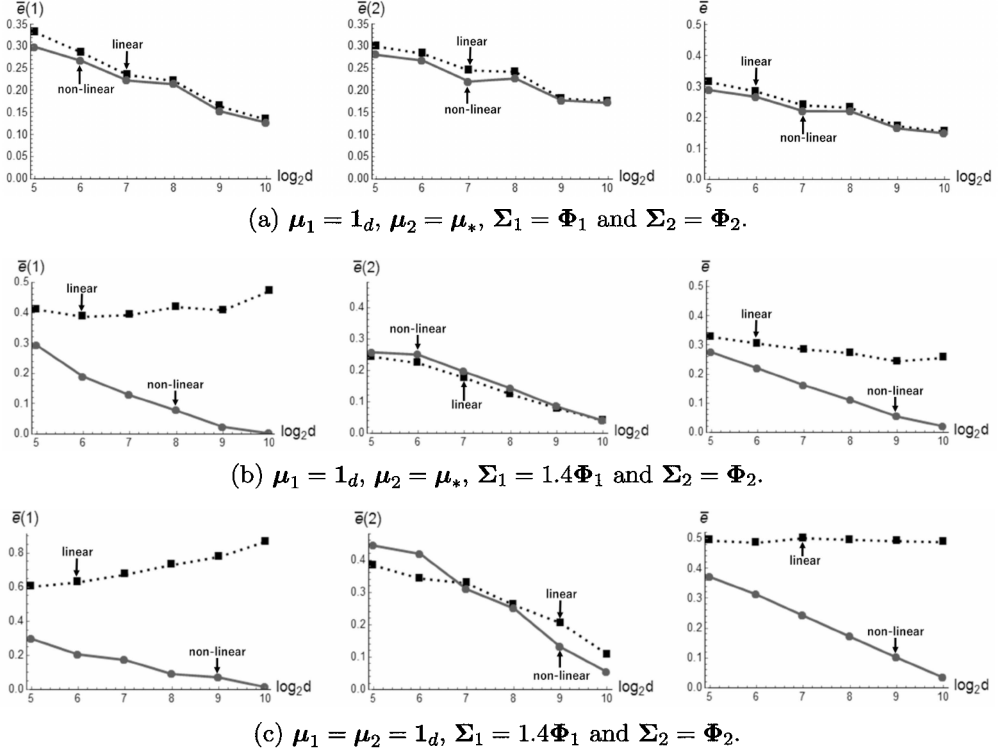


Figure 1: The performance of the linear SVM and the non-linear SVM for (a) to (c). The error rates of the linear SVM are denoted by the dotted lines, and those of the non-linear SVM are denoted by the solid lines.

4 Simulation

In this section, we compare the performance of the linear SVM given by (4) and the non-linear SVM given by (6) in numerical simulations.

We set $d = 2^s$, $s = 5, \dots, 10$, and $(n_1, n_2) = (5, 5)$. We generated \mathbf{x}_{ij} , $j = 1, 2, \dots$, ($i = 1, 2$) independently from $\pi_i : N_d(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$. We set $\boldsymbol{\mu}_* = (1, \dots, 1, 0, \dots, 0)^T$ whose last $\lceil d^{2/3} \rceil$ elements are 0, where $\lceil x \rceil$ denotes the smallest integer $\geq x$. Let $\boldsymbol{\Phi}_1 = \mathbf{B}(0.3^{|i-j|^{1/3}})\mathbf{B}$, $\boldsymbol{\Phi}_2 = \mathbf{B}(0.4^{|i-j|^{1/3}})\mathbf{B}$ and

$$\mathbf{B} = \text{diag}[\{0.5 + 1/(d+1)\}^{1/2}, \dots, \{0.5 + d/(d+1)\}^{1/2}].$$

We considered three cases :

- (a) $\boldsymbol{\mu}_1 = \mathbf{1}_d = (1, \dots, 1)^T$, $\boldsymbol{\mu}_2 = \boldsymbol{\mu}_*$, $\boldsymbol{\Sigma}_1 = \boldsymbol{\Phi}_1$ and $\boldsymbol{\Sigma}_2 = \boldsymbol{\Phi}_2$;
- (b) $\boldsymbol{\mu}_1 = \mathbf{1}_d$, $\boldsymbol{\mu}_2 = \boldsymbol{\mu}_*$, $\boldsymbol{\Sigma}_1 = 1.4\boldsymbol{\Phi}_1$ and $\boldsymbol{\Sigma}_2 = \boldsymbol{\Phi}_2$; and
- (c) $\boldsymbol{\mu}_1 = \boldsymbol{\mu}_2 = \mathbf{1}_d$, $\boldsymbol{\Sigma}_1 = 1.4\boldsymbol{\Phi}_1$ and $\boldsymbol{\Sigma}_2 = \boldsymbol{\Phi}_2$.

For $\mathbf{x}_0 \in \pi_i$ ($i = 1, 2$) we calculated each classifier 2000 times to confirm if each rule does (or does not) classify \mathbf{x}_0 correctly and defined $P_{ir} = 0$ (or 1) accordingly for each π_i . We calculated the error rates, $\bar{e}(i) = \sum_{r=1}^{2000} P_{ir}/2000$, $i = 1, 2$. Also, we calculated the average error rate, $\bar{e} = \{\bar{e}(1) + \bar{e}(2)\}/2$. For the Gaussian kernel, we chose γ from the candidates, $d^{(t+5)/10}$, $t = 1, \dots, 10$, by a cross-validation procedure. Their standard deviations are less than 0.011. In Figure 1, we plotted $\bar{e}(1)$, $\bar{e}(2)$ and \bar{e} for (a) to (c).

We observed that the SVMs give preferable performances for (a) in Figure 1. However, the linear SVM gave a quite bad performance for (c). This is because of $\Delta = 0$ for (c). On the other hand, the non-linear SVM gave a better performance compared to the linear SVM for (b) and (c). This is because the non-linear SVM draws information about heteroscedasticity from the difference of $\text{tr}(\boldsymbol{\Sigma}_i)$ s. See Section 3 for the details.

5 Appendix

Proof of Lemma 3.1. Similarly to the proof of Lemma 1 in Nakayama et al. [8], we have that as $d \rightarrow \infty$

$$L_*(\boldsymbol{\alpha}) = 2\alpha_* - \frac{\Delta_*}{2}\alpha_*^2\{1 + o_p(1)\} - \frac{1}{2}\left((1 - c_1)\sum_{j=1}^{n_1}\alpha_j^2 + (1 - c_2)\sum_{j=n_1+1}^N\alpha_j^2\right)$$

subject to (3) under (A-iv) to (A-vi), where $\alpha_* = \sum_{j=1}^{n_1}\alpha_j$. Then, by noting

$$\liminf_{d \rightarrow \infty} (1 - c_i)/\Delta_* > 0, \quad i = 1, 2,$$

under (A-iv), in a way similar to the proof of Lemma 2 in Nakayama et al. [8], we can obtain the result. \square

Proof of Theorem 3.1. By using (7), the result is obtained straightforwardly. \square

Acknowledgements

Research of the second author was partially supported by Grant-in-Aid for Young Scientists (B), Japan Society for the Promotion of Science (JSPS), under Contract Number 26800078. Research of the third author was partially supported by Grants-in-Aid for Scientific Research (A), JSPS, under Contract Number 15H01678.

References

- [1] Aoshima, M., Yata, K., 2011. Two-stage procedures for high-dimensional data. *Sequential Anal.* (Editor's special invited paper) 30, 356-399.
- [2] Aoshima, M., Yata, K., 2014. A distance-based, misclassification rate adjusted classifier for multiclass, high-dimensional data. *Ann. Inst. Statist. Math.* 66, 983-1010.
- [3] Aoshima, M., Yata, K., 2015a. Geometric classifier for multiclass, high-dimensional data. *Sequential Anal.* 34, 279-294.
- [4] Aoshima, M., Yata, K., 2015b. High-dimensional quadratic classifiers in non-sparse settings. *arXiv:1503.04549*.
- [5] Chan, Y.-B., Hall, P., 2009. Scale adjustments for classifiers in high-dimensional, low sample size settings. *Biometrika* 96, 469-478.
- [6] Hall, P., Marron, J.S., Neeman, A., 2005. Geometric representation of high dimension, low sample size data. *J. R. Statist. Soc. B* 67, 427-444.
- [7] Marron, J.S., Todd, M.J., Ahn, J., 2007. Distance-weighted discrimination. *J. Amer. Statist. Assoc.* 102, 1267-1271.
- [8] Nakayama, Y., Yata, K., Aoshima, M., 2017. Support vector machine and its bias correction in high-dimension, low-sample-size settings. Revised in *J. Stat. Plan. Infer.* (*arXiv:1702.08019*).
- [9] Schölkopf, B., Smola, A.J., 2002. *Learning with Kernels*. MIT Press, Cambridge.
- [10] Vapnik, V.N., 2000. *The Nature of Statistical Learning Theory* (second ed.). Springer, New York.